# INFO8002

# Large-Scale Data Systems

## Exercise Session #2

LIÈGE université

# Reminder

# REMINDER :

## Distributed System Models

Distributed system models are defined as the combination of **(i)** a <u>process abstraction</u>, **(ii)** a <u>link abstraction</u>, and **(iii)** a <u>failure detector abstraction</u>:

- **Fail-stop**
  - <u>Process</u> : **Crash-stop**
  - <u>Link</u> : **Perfect**
  - <u>Failure Detector</u> : **Perfect**

- **Fail-silent**
  - <u>Process</u> : **Crash-stop**
  - <u>Link</u> : **Perfect**
  - <u>Failure Detector</u> : **/**

- **Fail-noisy**
  - <u>Process</u> : **Crash-stop**
  - <u>Link</u> : **Perfect links**
  - <u>Failure Detector</u> : **Eventually perfect**

- **Fail-recovery**
  - <u>Process</u> : **Crash-recovery**
  - <u>Link</u> : **Stubborn**
  - <u>Failure Detector</u> : **/**

- **Fail-arbitrary**
  - <u>Process</u> : **Byzantine**
  - <u>Link</u> : **Perfect**
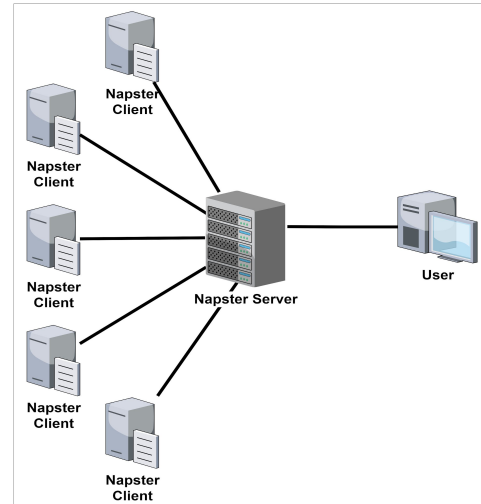  - <u>Failure Detector</u> : **/**

# REMINDER :

## Shared Memory Abstraction

- **Shared memory** can be viewed as **an array of registers** to which a process can read or write.

- **Shared memory models** are defined as a combination of behaviour of registers in presence of **(i)** failures, and **(ii)** concurrent operations :

  - **Safe Register (not seen)**
    - Failures : **??**
    - Concurrency : **Arbitrary** value

  - **Regular Register**
    - Failures : **Fail-stop or Fail-silent (others?)**
    - Concurrency: **Previous** or **concurrently** written value

  - **Atomic Register**
    - Failures : **Fail-stop (others?)**
    - Concurrency : Ensure **linearisability** of operations

# PROBLEM 1
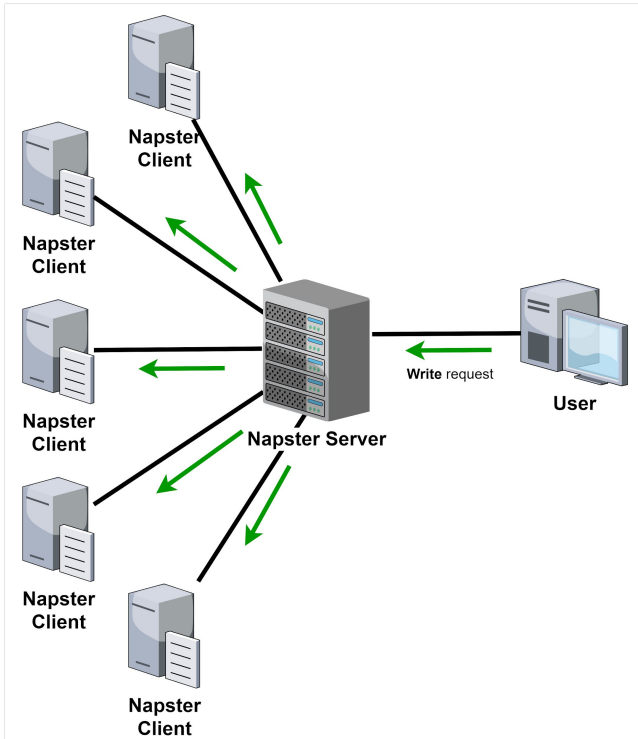
# Content Sharing System

You are responsible for designing a system allowing the **storage** and **distribution** of content.
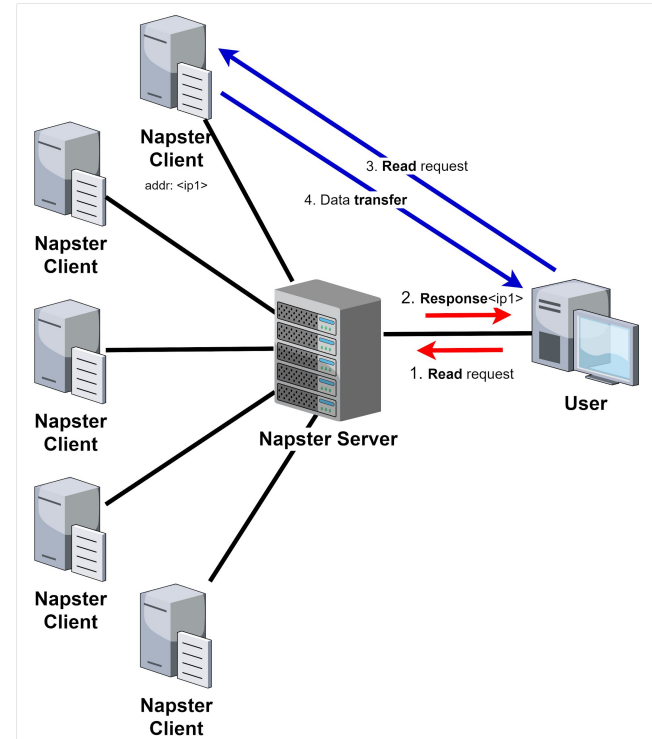


**Specify** an architecture for this distributed storage system and **provide** a pseudo-implementation using shared registers.
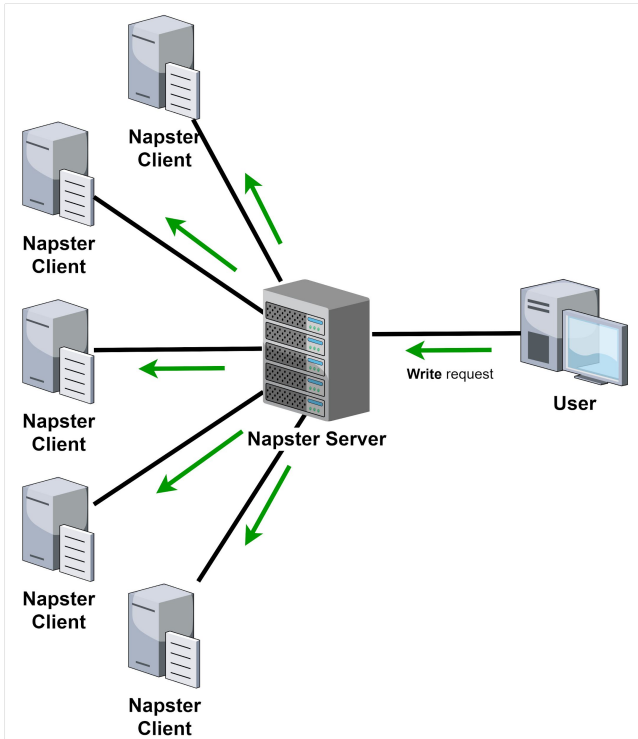
# Content Sharing System

## Write query



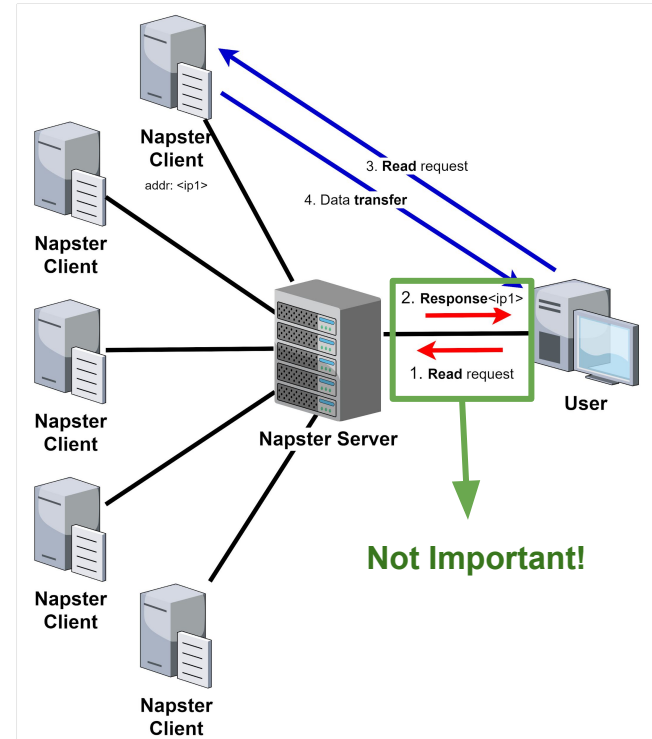**Napster Client**

**Napster Client**

**Napster Client**

**Napster Client**

**Napster Client**

**Napster Server**

**Write** request

**User**

## Read query



**Napster Client**

addr: <ip1>

**Napster Client**

**Napster Client**

**Napster Client**

**Napster Server**

3. **Read** request

4. Data **transfer**

2. **Response**<ip1>

1. **Read** request

**User**

# Content Sharing System

## Write query



**Napster Client** (×6)

**Napster Server**

**Write** request

**User**

## Read query



**Napster Client**

addr: <ip1>

3. **Read** request

4. Data **transfer**

**Napster Client** (×5)

**Napster Server**

2. **Response**<ip1>

1. **Read** request

**User**

**Not Important!**

# Content Sharing System

**Problems:**

1. How can we ensure that operations terminates ?

# Content Sharing System

**Problems:**

1. **How can we ensure that operations terminates ?**

   Property n°1: **"Termination"**

   *"If a correct process invokes an operation, then the operation eventually completes."*

2. **How can we ensure that users receives a <u>coherent</u> response to their read request ?**

# Content Sharing System

**Problems:**

1. **How can we ensure that operations terminates ?**

   Property n°1: **"Termination"**

   *"If a correct process invokes an operation, then the operation eventually completes."*

2. **How can we ensure that users receives a <u>coherent</u> response to their read request ?**

   Property n°2: **"Validity"**

   *"A **read** that is <u>not concurrent</u> with a **write** <u>returns the last value written</u>; a **read** that is <u>concurrent</u> with a **write** <u>returns the last value written or the value concurrently written</u>."*

# Content Sharing System

## Module Specification:

**Module 1:** Interface and properties of distributed storage

    **Module:**

        **Name:** *NapsterClientServer*, **instance** *np*.

    **Events**:

        **Request:** < np, Read | r, m > **:** <u>Invokes</u> a **read** operation on **m** consecutive registers starting on register **r**.

        **Request:** < np, Write | v, r> **:** <u>Invokes</u> a **write** operation with value **v** starting on register **r**.

        **Indication:** < np, ReadReturn | v > **:** <u>Completes</u> a **read** operation with return value **v**.

        **Indication:** < np, WriteReturn> **:** <u>Completes</u> a **write** operation.

    **Properties:**

        **NP1:** Termination.

        **NP2:** Validity.

*"(1, N) Regular Register"* <u>Module</u>

# Content Sharing System

**Implementation**

**Algorithm 1:**

>**Implements**:
>>*NapsterClientServer*, **instance** *np*.
>
>**Uses**:
>>(1, N)-RegularRegister, **instance** *onrr*.

**upon event** < *np*, *Init* > **do**
>*???*

**upon event** < *np*, Write | v, r > **do**
>*???*

**upon event** < *np*, Read | r, m > **do**
>*???*

**upon event** < *onrr*, ReadReturn | r, v> **do**
>*???*

**upon event** < *onrr*, WriteReturn | r > **do**
>*???*

# Content Sharing System

**Implementation**

**Algorithm 1:**

    **Implements**:

        *NapsterClientServer*, **instance** *np*.

    **Uses**:

        (1, N)-RegularRegister, **instance** *onrr*.

    **upon event** < *np*, *Init* > **do**

        *memory* := $[0]^{MemorySize}$;

        *pendingR* := ∅**;**

        *pendingWr* := ∅**;**

    **upon event** < *np*, Write | v, r > **do**

        **forall** v' ∈ v **do**

            *pendingWr* := *pendingWr* ∪ {r + **index**(v)};

        **forall** v' ∈ v **do**

            **trigger** < *onrr*, Write | v', r + **index**(v)>;

**upon event** < *np*, Read | r, m > **do**

    *ReadRet* := $[0]^{m}$;

    *offset* := r;

    **for** *i* **in range(m) do** *pendingR* := *pendingR* ∪ {r + i};

    **for** *i* **in range(m) do** **trigger** < *onrr, Read | r + i >;*

**upon event** < *onrr*, ReadReturn | r, v> **do**

    *pendingR := pendingR \ {r};*

    *ReadRet[r-offset] := v;*

    **if** *pendingR* ⊆ ∅ **then**

        **trigger** < *np*, ReadReturn | *ReadRet* >;

**upon event** < *onrr*, WriteReturn | r > **do**

    *pendingWr := pendingWr \ {r};*

    **if** *pendingWr* ⊆ ∅ **then**

        **trigger** < *np*, WriteReturn >;

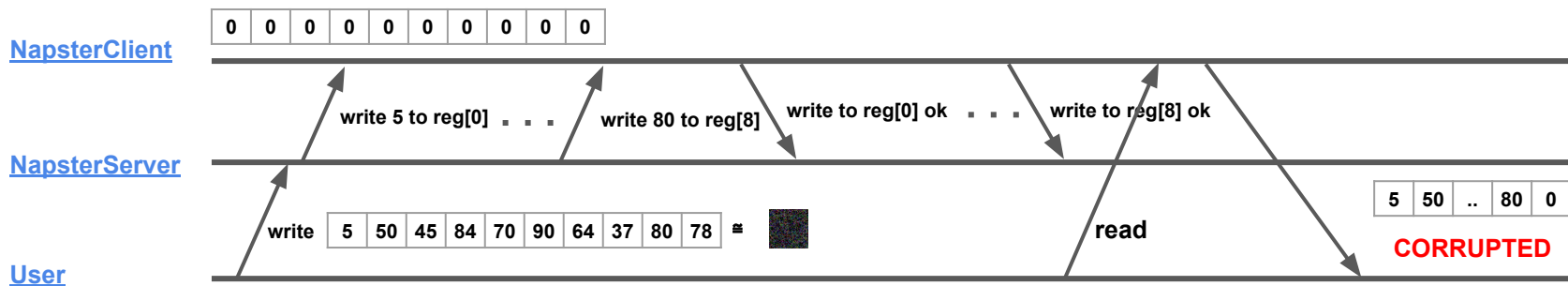# Content Sharing System

**PROBLEM:**



→ We would like to ensure **coherency** of the distributed storage.

# Content Sharing System

**PROBLEM:**



$\longrightarrow$ We would like to ensure **coherency** of the distributed storage.

**SOLUTION:** Write to **temporary buffer** until all data are received.

# Content Sharing System

**Implementation**

**Algorithm 2:**
    **Implements**:
        (1, N)-NewRegularRegister, **instance** *onnrr*.
    **Uses**:
        BestEffortBroadcast, **instance** *beb*;
        PerfectPointToPointLinks, **instance** *pl*;
        PerfectFailureDetector, **instance** *P*.

      **upon event** < *onnrr*, *Init* > **do**
            *val* := $[0]^{MemorySize}$;
            wrBuffer := $[0]^{MemorySize}$;
            *correct* := $\Pi$*;*
            *writeset* := ∅*;*

      **upon event** < *onnrr*, Write | v, r > **do**
            **trigger** < *beb*, Broadcast | [WRITE, v, r]>;

**upon event** < *beb*, Deliver | q, [Write, v, r]> **do**
    wrBuffer[r] = v;
    **trigger** < *pl*, Send | q, ACK>;

**upon event** < *onnrr*, Flush > **do**
    val := wrBuffer;
    **trigger** < *onnrr*, FlushReturn >;

**… (cfr. Theoretical Lectures : "*(1,N)-RegularRegister*")**

# Content Sharing System

**Implementation**

---

**Algorithm 3:**

    **Implements**:

        *NapsterClientServer*, **instance** *np*.

    **Uses**:

        (1, N)-NewRegularRegister, **instance** *onnrr*.

**upon event** $< np$, *Init* $>$ **do**

    *pendingR* **:=** *pendingWr* **:=** ∅**;**

**upon event** $< np$, Write $\mid$ v, r $>$ **do**

    **forall** v' ∈ v **do**

        *pendingWr* := *pendingWr* ∪ {r + **index**(v)};

    **forall** v' ∈ v **do**

        **trigger** $<$ *onnrr*, Write $\mid$ v', r + **index**(v)$>$;

**upon event** $< np$, Read $\mid$ r, m $>$ **do**

    *ReadRet* := $[0]^{m}$; *offset* := *r*;

    **for** i **in range(m) do** *pendingR* := *pendingR* ∪ {r + i};

    **for** i **in range(m) do trigger** $<$ *onrr*, Read $\mid r + i >$;

**upon event** $<$ *onrr*, ReadReturn $\mid$ r, v$>$ **do**

    *pendingR := pendingR \ {r};*

    *ReadRet[r-offset] := v;*

    **if** *pendingR* ⊆ ∅ **then**

        **trigger** $< np$, ReadReturn $\mid$ *ReadRet* $>$;

**upon event** $<$ *onrr*, WriteReturn $\mid$ r $>$ **do**

    *pendingWr := pendingWr \ {r};*

    **if** *pendingWr* ⊆ ∅ **then**

        **trigger** $<$ *onnrr*, Flush$>$

**upon event** $<$ *onnrr*, FlushReturn $>$ **do**

    **trigger** $< np$, WriteReturn $>$;

# Content Sharing System

**Problems:**

- **Do you see any other problems?**

# Content Sharing System

**Problems:**

- **Napster** is a technology allowing to have several servers that can be reached for write **and** read requests. How would you implement this in:
    1. In a **fail-stop** system?
    2. In a **fail-silent** system?
    3. In a **byzantine** system?

# HOMEWORK !